

Lecture 13 Lazy Training

*Instructor: Yingyu Liang**Date:**Scriber: Yuchen Zeng*

1 Lazy Training of General Neural Networks

In previous lectures, we have shown many good properties of two-layer neural networks using the NTK formulation. In this lecture, we will show the key technical lemma still holds for more general neural networks, and the intuition is still the same. The formulation here is usually called lazy training.

Consider n samples $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ and $\mathbf{y} = (y_1, \dots, y_n)^\top$. Denote the model of the neural network as function $f(\mathbf{x}, \mathbf{w})$ for weight $\mathbf{w} \in \mathbb{R}^p$, and let

$$\mathbf{f}(\mathbf{w}) = (f(\mathbf{x}_1, \mathbf{w}), \dots, f(\mathbf{x}_n, \mathbf{w}))^\top \in \mathbb{R}^n.$$

Note that the loss function

$$L(\alpha \mathbf{f}(\mathbf{w})) = \frac{1}{2} \|\alpha \mathbf{f}(\mathbf{w}) - \mathbf{y}\|^2,$$

and denote $L_0 = L(\alpha \mathbf{f}(\mathbf{w}(0)))$. Besides, we have

$$\frac{d\mathbf{w}(t)}{dt} = \dot{\mathbf{w}}(t) = -\nabla_{\mathbf{w}} L(\alpha \mathbf{f}(\mathbf{w})) = -\alpha \mathbf{J}_t^\top (\nabla L)(\alpha \mathbf{f}(\mathbf{w}(t))),$$

where $\nabla_{\mathbf{w}}$ is w.r.t. \mathbf{w} , ∇L is w.r.t. z in $L(z)$, and the Jacobian matrix \mathbf{J}_t is related to the Hessian matrix introduced in our previous lectures (more precisely, $\mathbf{J}_t \mathbf{J}_t^\top$ is equivalent to the Hessian). To be more specific, we have

$$\mathbf{J}_t = \mathbf{J}_{\mathbf{w}(t)} = (\nabla f(\mathbf{x}_1, \mathbf{w}(t)), \dots, \nabla f(\mathbf{x}_n, \mathbf{w}(t)))^\top \in \mathbb{R}^{n \times p}.$$

We will also introduce another gradient flow, which is on a linear approximation of the network. By the first-order Taylor expansion, we obtain

$$\mathbf{f}_0(\mathbf{u}) = \mathbf{f}(\mathbf{w}(0)) + \mathbf{J}_0(\mathbf{u} - \mathbf{w}(0)), \text{ where } \mathbf{u}(0) = \mathbf{w}(0).$$

Now we introduce the gradient flow for training the linear function $\mathbf{f}_0(\mathbf{u})$:

$$\frac{d\mathbf{u}(t)}{dt} = \dot{\mathbf{u}}(t) = -\nabla_{\mathbf{u}} L(\alpha \mathbf{f}_0(\mathbf{u}(t))) = -\alpha \mathbf{J}_0^\top (\nabla L)(\alpha \mathbf{f}_0(\mathbf{u}(t))).$$

The goal of this lecture is to show that the loss decreases exponentially fast. Remember that in our last lecture, if we do random initialization, we showed that the spectrum of the Hessian matrix is far away from 0, with high probability. To generalize the results to general neural networks, we first introduce some assumption.

Assumption 1.

$$\begin{aligned} \text{rank}(\mathbf{J}_0) &= n, p > n, \\ \sigma_{\min} &\triangleq \sigma_{\min}(\mathbf{J}_0) = \sqrt{\lambda_{\min}(\mathbf{J}_0\mathbf{J}_0^\top)} > 0, \\ \sigma_{\max} &\triangleq \sigma_{\max}(\mathbf{J}_0). \\ \exists \beta > 0, &\|\mathbf{J}_{\mathbf{w}} - \mathbf{J}_{\mathbf{v}}\| \leq \beta\|\mathbf{w} - \mathbf{v}\|. \end{aligned}$$

Next, we aim to show the following theorem.

Theorem 2. Under Assumption 1, if the scaling parameter $\alpha \geq \frac{6\beta\sqrt{8\sigma_{\max}^2 L_0}}{\sigma_{\min}^3}$, we have

$$\begin{aligned} \max\{L(\alpha\mathbf{f}(\mathbf{w}(t))), L(\alpha\mathbf{f}_0(\mathbf{w}(t)))\} &\leq L_0 \exp\left(-t\frac{\alpha^2\sigma_{\min}^2}{2}\right), \\ \max\{\|\mathbf{w}(t) - \mathbf{w}(0)\|, \|\mathbf{u}(t) - \mathbf{u}(0)\|\} &\leq \frac{3\sqrt{8\sigma_{\max}^2 L_0}}{\alpha\sigma_{\min}^2}. \end{aligned}$$

To prove Theorem 2, we introduce the following lemmas.

Lemma 3.

$$\begin{aligned} \frac{d}{dt}\alpha\mathbf{f}(\mathbf{w}(t)) &= \alpha\mathbf{J}_t\dot{\mathbf{w}}(t) = -\alpha^2\mathbf{J}_t\mathbf{J}_t^\top\nabla L(\alpha\mathbf{f}(\mathbf{w}(t))) \\ &= -\alpha^2\mathbf{J}_t\mathbf{J}_t^\top(\alpha\mathbf{f}(\mathbf{w}(t)) - \mathbf{y}) \\ \frac{d}{dt}\alpha\mathbf{f}_0(\mathbf{u}(t)) &= -\alpha\mathbf{J}_0\mathbf{J}_0^\top(\alpha\mathbf{f}_0(\mathbf{u}(t)) - \mathbf{y}). \end{aligned}$$

We consider some general dynamics as follows. Once we consider this more general setting, we can see the intuition more clearly.

Lemma 4. Suppose $\dot{\mathbf{z}}(t) = -Q(t)\nabla L(\mathbf{z}(t)), \forall t \in [0, T]$. If $\lambda = \inf_{t \in [0, T]} \lambda_{\min}(Q_t) > 0$, then for $t \in [0, T]$, $L(\mathbf{z}(t)) \leq L(\mathbf{z}(0)) \exp(-2\lambda t)$.

Proof.

$$\begin{aligned} \frac{d}{dt}L(\mathbf{z}(t)) &= \frac{d}{dt}\frac{1}{2}\|\mathbf{z}(t) - \mathbf{y}\|^2 \\ &= \frac{1}{2}\langle \dot{\mathbf{z}}(t), \mathbf{z}(t) - \mathbf{y} \rangle \\ &= \langle -Q_t(\mathbf{z}(t) - \mathbf{y}), \mathbf{z}(t) - \mathbf{y} \rangle \\ &\leq -\lambda_{\min}(Q_t)\|\mathbf{z}(t) - \mathbf{y}\|^2 \\ &\leq -\lambda\|\mathbf{z}(t) - \mathbf{y}\|^2 \\ &= 2\lambda L(\mathbf{z}(t)). \end{aligned}$$

Next, we can use Grönwall's inequality to get the final bound. □

Lemma 5. Suppose $\dot{\mathbf{v}}(t) = -\mathbf{s}(t)^\top \nabla L(g(\mathbf{v}(t)))$ where $\mathbf{s}(t) = \mathbf{S}_t$ which is the Jacobian of g . Let $Q_t = \mathbf{S}_t \mathbf{S}_t^\top$ and assume $\forall t \in [0, T], \lambda_i(Q_t) \in [\lambda, \lambda_1]$. Then $\forall t \in [0, T]$,

$$\|\mathbf{v}(t) - \mathbf{v}(0)\|^2 \leq \frac{2\lambda_1 L(g(\mathbf{v}(0)))}{\lambda}.$$

Proof.

$$\begin{aligned} \|\mathbf{v}(t) - \mathbf{v}(0)\| &= \left\| \int_0^t \dot{\mathbf{v}}(s) ds \right\| \\ &\leq \int_0^t \|\dot{\mathbf{v}}(s)\| ds \\ &= \int_0^t \|\mathbf{s}(s)^\top (g(\mathbf{v}(s)) - \mathbf{y})\| ds \\ &\leq \sqrt{\lambda_1} \int_0^t \|g(\mathbf{v}(s)) - \mathbf{y}\| ds \\ &\leq \sqrt{\lambda_1} \int_0^t \|g(\mathbf{v}(0)) - \mathbf{y}\| \exp(-\lambda s) ds \end{aligned}$$

where the last step is by applying Lemma 4 on $\mathbf{z}(t) = g(\mathbf{v}(t))$ and noting $\dot{\mathbf{z}}(t) = -\mathbf{S}_t \mathbf{S}_t^\top \nabla L(\mathbf{z}(t))$. \square

By the two lemmas above, we can easily get the results for the linear function $\mathbf{f}_0(\mathbf{u})$:

$$\begin{aligned} \frac{d}{dt} \alpha \mathbf{f}_0(\mathbf{u}(t)) &= -\alpha^2 \mathbf{J}_0 \mathbf{J}_0^\top (\alpha \mathbf{f}_0(\mathbf{u}(t)) - \mathbf{y}) \\ L(\alpha \mathbf{f}_0(\mathbf{u}(t))) &\leq L_0 \exp(-2t\alpha^2 \sigma_{\min}^2) \\ \|\mathbf{u}(t) - \mathbf{u}(0)\| &\leq \frac{\sqrt{2\alpha^2 \sigma_{\max}^2 L_0}}{\alpha^2 \sigma_{\min}}. \end{aligned}$$

To get the results for $\mathbf{f}(\mathbf{w})$, we will need to show that \mathbf{J}_t is bounded.

Lemma 6. If \mathbf{w} satisfies $\|\mathbf{w} - \mathbf{w}(0)\| \leq B \triangleq \sigma_{\min}/(2\beta)$, then

$$\sigma_{\min}(\mathbf{J}_w) \geq \frac{\sigma_{\min}}{2}, \quad \sigma_{\max}(\mathbf{J}_w) \leq \frac{3\sigma_{\max}}{2}.$$

Proof. First consider $\sigma_{\max}(\mathbf{J}_w) = \|\mathbf{J}_w\|$. Recall that β is the Lipschitz constant for the Jacobian.

$$\begin{aligned} \|\mathbf{J}_w\| &= \|\mathbf{J}_w - \mathbf{J}_0 + \mathbf{J}_0\| \\ &\leq \|\mathbf{J}_w - \mathbf{J}_0\| + \|\mathbf{J}_0\| \\ &= \beta B + \sigma_{\max} \\ &\leq \frac{\sigma_{\min}}{2} + \sigma_{\max} \\ &\leq \frac{3\sigma_{\max}}{2}. \end{aligned}$$

Next consider $\sigma_{\min}(\mathbf{J}_w)$. By definition,

$$\sigma_{\min}^2(\mathbf{J}_w) = \lambda_{\min}(\mathbf{J}_w \mathbf{J}_w^\top) = \min_{\|\mathbf{v}\|=1} \mathbf{v}^\top \mathbf{J}_w \mathbf{J}_w^\top \mathbf{v} = \min_{\|\mathbf{v}\|=1} \|\mathbf{J}_w^\top \mathbf{v}\|^2.$$

Let $\mathbf{A}_v = (\mathbf{J}_w - \mathbf{J}_0)^\top \mathbf{v}$, and $\mathbf{B}_v = \mathbf{J}_0^\top \mathbf{v}$. Then we have

$$\begin{aligned} \|\mathbf{J}_w^\top \mathbf{v}\|^2 &= \|(\mathbf{J}_w - \mathbf{J}_0 + \mathbf{J}_0)^\top \mathbf{v}\|^2 \\ &= \|\mathbf{A}_v + \mathbf{B}_v\|^2 \\ &= \|\mathbf{A}_v\|^2 + 2\langle \mathbf{A}_v, \mathbf{B}_v \rangle + \|\mathbf{B}_v\|^2 \\ &\geq \|\mathbf{A}_v\|^2 - 2\|\mathbf{A}_v\|\|\mathbf{B}_v\| + \|\mathbf{B}_v\|^2 \\ &= (\|\mathbf{A}_v\| - \|\mathbf{B}_v\|)^2. \end{aligned} \tag{1}$$

Note that

$$\begin{aligned} \|\mathbf{B}_v\| &= \|\mathbf{J}_0^\top \mathbf{v}\| \geq \sigma_{\min}, \\ \|\mathbf{A}_v\| &\leq \|(\mathbf{J}_w - \mathbf{J}_0)^\top \mathbf{v}\| \leq \|\mathbf{J}_w - \mathbf{J}_0\| \|\mathbf{v}\| = \|\mathbf{J}_w - \mathbf{J}_0\| \leq \beta B = \sigma_{\min}/2. \end{aligned} \tag{2}$$

Combining equation 1 and equation 2 leads to the desired results. \square

With this lemma which bounds the spectrum and the other two key lemmas above, we have

$$\begin{aligned} \frac{d}{dt} \alpha \mathbf{f}(\mathbf{w}(t)) &= -\alpha^2 \mathbf{J}_t \mathbf{J}_t^\top (\alpha \mathbf{f}(\mathbf{w}(t)) - \mathbf{y}) \\ \lambda &= \alpha^2 \frac{\sigma_{\min}^2}{4} \\ L(\alpha \mathbf{f}(\mathbf{w}(t))) &\leq L_0 \exp(-t\alpha^2 \sigma_{\min}^2/2) \\ \|\mathbf{w}(t) - \mathbf{w}(0)\| &\leq \frac{\sqrt{\frac{9\sigma_{\min}^2}{2} \alpha^2 L_0}}{\alpha^2 \sigma_{\min}^2/4} = \frac{2\sqrt{8\sigma_{\max}^2 L_0}}{\alpha \sigma_{\min}^2} \triangleq B'. \end{aligned}$$

Now to complete the proof of the Theorem 2, it is sufficient to ensure $B' \leq B$. This is guaranteed by the condition on α in the theorem.